

# Technical Note for Discrete-Time Diffusion Approximations Motivated from Hospital Inpatient Flow Management

J. G. Dai<sup>1</sup> and Pengyi Shi<sup>2</sup>

<sup>1</sup>School of Operations Research and Information Engineering, Cornell University  
[jd694@cornell.edu](mailto:jd694@cornell.edu)

<sup>2</sup>Krannert School of Management, Purdue University  
[shi178@purdue.edu](mailto:shi178@purdue.edu)

## Abstract

This note details the development of a discrete-time diffusion process to approximate the midnight customer count process in a  $M_{\text{per}}/\text{Geo}_{2\text{timeScale}}/N$  system. We prove a limit theorem that supports this diffusion approximation, and discuss two methods to compute the stationary distribution of this discrete-time diffusion process.

The  $M_{\text{per}}/\text{Geo}_{2\text{timeScale}}/N$  is single-pool queueing system with a periodic Poisson arrival process and a *two-time-scale* service time feature. This queueing system is motivated to study hospital inpatient flow management and is introduced in [6] (which we will refer to as the “main paper”). To analyze this system, a critical step is to obtain the stationary distribution  $\pi$  for the midnight count process  $\{X_k, k = 0, 1, \dots\}$ , where  $X_k$  denotes the number of customers in the system at the midnight of day  $k$ , including both customers in service and those waiting in the buffer. To efficiently compute  $\pi$ , especially when the number of servers  $N$  is large and the utilization is close to 1, we develop a discrete-time diffusion process  $\{X_k^*, k = 0, 1, \dots\}$  to approximate the midnight count process, and use the stationary distribution of  $\{X_k^*\}$  to approximate  $\pi$ .

In this note, we first prove a limit theorem in Section 1 that supports this diffusion approximation. Then, in Section 2 we discuss two methods to numerically compute or approximate the stationary distribution of the discrete-time diffusion process  $\{X_k^*\}$ . Finally, in Section 3, we show the accuracy of these two methods in approximating  $\pi$  via numerical experiments.

## 1 Diffusion limits for the single-pool model

Section 4.3 of [6] proposes a discrete-time diffusion process to approximate the midnight count process. This approximation is motivated by a limit theorem that shows the convergence of stochastic processes. In this section, we prove this limit theorem.

Instead of fixing the number of servers  $N$ , we consider a sequence of  $M_{\text{per}}/\text{Geo}_{2\text{timeScale}}/N$  systems indexed by  $N$ , i.e., a sequence of the single-pool models described in the main paper [6]. Let  $\Lambda^N$  be the daily arrival rate of the  $N$ th system. Let  $m = 1/\mu$ , the mean LOS, be fixed and  $\rho^N = (\Lambda^N m)/N$  be the traffic intensity of the  $N$ th system. We assume that

$$\lim_{N \rightarrow \infty} \Lambda^N/N = \Lambda^*, \text{ and } \lim_{N \rightarrow \infty} \sqrt{N}(1 - \rho^N) = \beta^* \text{ for some } \beta^* > 0. \quad (1)$$

Analogous to the conventional many-server queues that model customer call centers [7], we call Condition (1) the *Quality- and Efficiency-Driven* (QED) condition.

We use  $X_k^N$  to denote the midnight customer count at the midnight of day  $k$  in the  $N$ th system. We consider the *diffusion-scaled* midnight customer count processes  $\tilde{X}^N = \{\tilde{X}_k^N : k = 0, 1, 2, \dots\}$

for the sequence of single-pool systems, where for a given  $k$ ,  $\tilde{X}_k^N$  is defined as

$$\tilde{X}_k^N = \frac{X_k^N - N}{\sqrt{N}}. \quad (2)$$

Adapting the derivations in the main paper, we can show that  $\tilde{X}_k^N$  satisfies the following relationship:

$$\tilde{X}_k^N = \tilde{Y}_k^N + \mu \sum_{i=0}^{k-1} (\tilde{X}_i^N)^-, \quad k = 0, 1, 2, \dots, \quad (3)$$

where

$$\tilde{Y}_k^N = \tilde{X}_0^N + \frac{1}{\sqrt{N}} \left( A_{(0,k]}^N - k\Lambda^N \right) - \frac{1}{\sqrt{N}} \left( D_{(0,k]}^N - \mu(Z_0^N + \dots + Z_{k-1}^N) \right) + k\sqrt{N}\mu(\rho^N - 1),$$

$A_{(0,k]}^N = \sum_{i=0}^{k-1} A_i^N$  and  $D_{(0,k]}^N = \sum_{i=0}^{k-1} D_i^N$  are the cumulative number of arrivals and departures from 0 until the midnight (zero hour) of day  $k$  in the  $N$ th system, respectively, and  $Z_i^N = \min(X_i^N, N)$  is the number of busy servers at the midnight of day  $i$ . We assume the initial condition

$$\tilde{X}_0^N \Rightarrow X_0^* \text{ as } N \rightarrow \infty, \quad (4)$$

where  $\Rightarrow$  denotes convergence in distribution. Under the many-server heavy-traffic framework (e.g., see [4]), we prove the following limit theorem:

**Theorem 1** *Consider a sequence of  $M_{\text{peri}}/\text{Geo}_{2\text{timeScale}}/N$  single-pool systems that satisfies (1) and (4). For any positive integer  $K \in \mathbb{Z}_+$ ,  $\tilde{X}^N \Rightarrow X^\dagger$  on the compact set  $[0, K]$  as  $N \rightarrow \infty$ , i.e.,*

$$\left( \tilde{X}_0^N, \tilde{X}_1^N, \dots, \tilde{X}_K^N \right) \Rightarrow \left( X_0^\dagger, X_1^\dagger, \dots, X_K^\dagger \right) \text{ as } N \rightarrow \infty. \quad (5)$$

The discrete-time limit process  $X^\dagger = \{X_k^\dagger, k = 0, 1, \dots\}$  satisfies

$$X_k^\dagger = Y_k^\dagger + \mu \sum_{i=0}^{k-1} (X_i^\dagger)^-, \quad k = 0, 1, \dots, \quad (6)$$

where  $Y_k^\dagger = Y^\dagger(k)$  for  $k = 0, 1, \dots$ , is an embedding of the Brownian motion  $\{Y^\dagger(t), t \geq 0\}$  which starts from  $X_0^\dagger$  and has mean  $-\mu\beta$  and variance  $\Lambda^* + \mu(1 - \mu)$ .

In this limit theorem, we deliberately use the superscript  $\dagger$  to differentiate the limit process  $X^\dagger$  (and the associated process  $Y^\dagger$ ) from the diffusion approximation  $X^*$  (and the associated process  $Y^*$ ) introduced in Section 4.3 of the main paper.

The key step of the proof for Theorem 1 is to show that  $\{\tilde{Y}_k^N, k = 0, 1, \dots\}$  converges to  $\{Y_k^\dagger, k = 0, 1, \dots\}$  on any given compact set  $[0, K]$ , or equivalently,

$$\left( \tilde{Y}_0^N, \tilde{Y}_1^N, \dots, \tilde{Y}_K^N \right) \Rightarrow \left( Y_0^\dagger, Y_1^\dagger, \dots, Y_K^\dagger \right) \text{ as } N \rightarrow \infty. \quad (7)$$

Then, the convergence of  $\tilde{X}^N$  to  $X^\dagger$  naturally follows because of the linear forms in (3) and (6). To prove (7), we first prove the convergence of the diffusion-scaled arrival processes in Section 1.1, and then the convergence of the discharge processes in Section 1.2.

## 1.1 Arrival process

For the  $N$ th system, let  $\tilde{E}_k^N = \frac{1}{\sqrt{N}} \left( A_{(0,k]}^N - k\Lambda^N \right)$ . We also introduce a continuous-time process  $\{\tilde{E}^N(t), t \geq 0\}$  defined as

$$\tilde{E}^N(t) = \frac{1}{\sqrt{N}} \left( E^N(t) - \Lambda^N t \right), \quad (8)$$

where  $E^N(\cdot)$  represents a Poisson process with rate  $\Lambda^N$ . It is easy to verify that  $\{\tilde{E}_k^N\}$  is an embedding of  $\tilde{E}^N(\cdot)$ , i.e.,

$$\tilde{E}_k^N = \tilde{E}^N(k), \quad k = 0, 1, \dots$$

Following a standard functional central limit theorem argument, we can show that

$$\tilde{E}^N(\cdot) \Rightarrow E^\dagger(\cdot) \quad (9)$$

in space  $\mathbb{D}$  endowed with the Skorohod  $J_1$  topology, where  $E^\dagger(\cdot)$  is a Brownian motion with drift 0 and variance  $\Lambda^*$ . Because the convergence of stochastic processes implies the convergence of any finite-dimensional joint distributions, we then naturally have

$$\left( \tilde{E}_0^N, \tilde{E}_1^N, \dots, \tilde{E}_K^N \right) \Rightarrow \left( E_0^\dagger, E_1^\dagger, \dots, E_K^\dagger \right) \text{ as } N \rightarrow \infty, \quad (10)$$

where  $E_k^\dagger = E^\dagger(k)$  is also an embedding of  $E^\dagger(\cdot)$ .

## 1.2 Discharge process

Now we consider the diffusion-scaled discharge processes. For the  $N$ th system, we introduce two discrete-time processes:

$$\check{D}_k^N = \frac{1}{\sqrt{N}} \left( D_{(0,k]}^N - \mu(Z_0^N + \dots + Z_{k-1}^N) \right), \quad k = 0, 1, \dots,$$

and

$$\tilde{D}_k^N = \frac{1}{\sqrt{N}} \sum_{i=1}^{Z_0^N + \dots + Z_{k-1}^N} (\xi_i - \mu), \quad k = 0, 1, \dots,$$

where  $\{\xi_i\}$  is a sequence of iid Bernoulli random variables with success probability  $\mu$ . Recall that in Appendix C of the main paper, we establish a revised system which tosses coins for every customer in service at the midnight to determine the departures each day, and we have proved this revised system is equivalent to the original system in distribution. Using the revised system, we can show the above two discrete-time processes are equal in distribution, i.e.,

$$(\check{D}_0^N, \check{D}_1^N, \dots) =^d (\tilde{D}_0^N, \tilde{D}_1^N, \dots).$$

Thus, it is sufficient to prove for any given  $K \in \mathbb{Z}_+$ ,

$$\left( \tilde{D}_0^N, \tilde{D}_1^N, \dots, \tilde{D}_K^N \right) \Rightarrow (S_0^*, S_1^*, \dots, S_K^*) \text{ as } N \rightarrow \infty. \quad (11)$$

Here,  $S_k^* = S^*(k)$  is an embedding of the Brownian motion  $S^*(\cdot)$  with drift 0 and variance  $\mu(1 - \mu)$ .

Let  $\eta_i = \xi_i - \mu$ , and  $\{\eta_i\}$  forms a sequence of iid random variables with mean 0 and variance  $\mu(1 - \mu)$ . We also define

$$T_k^N = \sum_{j=0}^{k-1} Z_j, \quad \bar{T}_k^N = \frac{T_k^N}{N},$$

and

$$S_n = \sum_{i=1}^n \eta_i.$$

Then, we can further rewrite  $\tilde{D}_k^N$  as

$$\tilde{D}_k^N = \frac{1}{\sqrt{N}} \sum_{i=1}^{T_k^N} \eta_i = \frac{1}{\sqrt{N}} S_{\bar{T}_k^N N}.$$

Correspondingly, proving (11) is equivalent to showing

$$\left( \frac{1}{\sqrt{N}} S_{\bar{T}_0^N N}, \frac{1}{\sqrt{N}} S_{\bar{T}_1^N N}, \dots, \frac{1}{\sqrt{N}} S_{\bar{T}_K^N N} \right) \Rightarrow (S_0^*, S_1^*, \dots, S_K^*) \text{ as } N \rightarrow \infty. \quad (12)$$

To prove (12), we introduce a continuous-time process  $\{\tilde{S}^N(t), t \geq 0\}$ , where

$$\tilde{S}^N(t) = \frac{1}{\sqrt{N}} S_{[tN]} \circ \bar{T}_{[t]}^N, \quad t \geq 0.$$

In other words,  $\tilde{S}^N(\cdot)$  is a composition of two continuous processes,  $\frac{1}{\sqrt{N}} S_{[\cdot N]}$  and  $\bar{T}_{[\cdot]}^N$ . It is easy to verify that  $\{\tilde{D}_k^N\}$  is an embedding of  $\tilde{S}^N(\cdot)$  because when  $t = k$ ,  $\bar{T}_k^N N = T_k^N$  is always an integer and

$$\tilde{D}_k^N = \frac{1}{\sqrt{N}} S_{\bar{T}_k^N N} = \tilde{S}^N(k).$$

If we can show

$$\frac{1}{\sqrt{N}} S_{[\cdot N]} \Rightarrow S^*(\cdot) \quad (13)$$

in space  $\mathbb{D}$  endowed with the Skorohod  $J_1$  topology as well as

$$\bar{T}_{[\cdot]}^N \rightarrow \bar{T}_{[\cdot]} \text{ in probability} \quad (14)$$

with  $\bar{T}_{[t]} = [t]$ , then applying the random time change theorem, we can prove (12). The convergence in (13) follows from the Donsker's theorem, and we focus on proving (14) below. It is sufficient to show for each  $0 \leq k \leq K$ ,  $Z_k^N/N \rightarrow 1$  almost surely, which we prove with induction.

We first rewrite the system equation under the fluid scaling:

$$\bar{X}_k^N = \bar{Y}_k^N + \sum_{i=0}^{k-1} (\bar{X}_i^N)^-, \quad (15)$$

where

$$\bar{X}_k^N = \frac{X_k^N - N}{N},$$

and

$$\bar{Y}_k^N = \bar{X}_0^N + \frac{\sum_{i=1}^k (A_{i-1}^N - \Lambda^N)}{N} - \frac{\sum_{i=1}^{T_k^N} \eta_i}{N} + \frac{k(\rho^N - 1)}{\sqrt{N}}.$$

Assume that  $X_0^N = N$ , then  $\bar{X}_0^N = 0$  and  $Z_0^N = N$  (so  $\bar{X}_0^N \rightarrow 0$  is trivial).

- When  $k = 1$ , we have

$$\bar{Y}_1^N = \bar{X}_0^N + \frac{A_0^N - \Lambda^N}{N} - \frac{\sum_{i=1}^N \eta_i}{N} + \frac{(\rho^N - 1)}{\sqrt{N}}.$$

Recall that  $(A_0^N - \Lambda^N)$  and  $\eta_i$  are centered random variables with mean 0. By the Law of Large Numbers, it is obvious that

$$\bar{Y}_1^N \rightarrow 0 \quad a.s. \text{ when } N \rightarrow \infty.$$

Thus,  $\bar{X}_1^N = \bar{Y}_1^N \rightarrow 0 \text{ a.s.}$ , and  $Z_1^N/N \rightarrow 1 \text{ a.s.}$ .

- Assume that at  $k$ , we have for all  $0 \leq j \leq k$ ,  $\bar{X}_j^N \rightarrow 0 \text{ a.s.}$  and  $Z_j^N/N \rightarrow 1 \text{ a.s.}$ . Then for  $k + 1$ , we have  $\bar{T}_{k+1}^N \rightarrow (k + 1) \text{ a.s.}$  and

$$\begin{aligned} \bar{Y}_{k+1}^N &= \bar{X}_0^N + \frac{\sum_{i=1}^{k+1} (A_{i-1}^N - \Lambda^N)}{N} - \frac{\sum_{i=1}^{T_{k+1}^N} \eta_i}{N} + \frac{(k+1)(\rho^N - 1)}{\sqrt{N}} \\ &= \frac{\sum_{i=1}^{k+1} (A_{i-1}^N - \Lambda^N)}{N} - \frac{\sum_{i=1}^{T_{k+1}^N} \eta_i}{T_{k+1}^N} \cdot \frac{T_{k+1}^N}{N} + \frac{(k+1)(\rho^N - 1)}{\sqrt{N}} \end{aligned} \quad (16)$$

$$\rightarrow 0 \quad a.s.. \quad (17)$$

Then

$$\bar{X}_{k+1}^N = \bar{Y}_{k+1}^N + \sum_{j=0}^k (\bar{X}_j^N)^- \rightarrow 0 \quad a.s.,$$

which completes the proof of Theorem 1.

## 2 Computing the stationary distribution of the discrete-time diffusion process

Motivated by the limit theorem proved in Section 1, Section 4.3 of the main paper [6] proposes a discrete-time diffusion process  $\{X_k^*, k = 0, 1, \dots\}$  to approximate the original midnight count process  $\{X_k, k = 0, 1, \dots\}$ . The dynamics of this approximation process follows:

$$X_k^* = Y_k^* + \mu \sum_{i=0}^{k-1} (X_i^*)^-, \quad k = 0, 1, 2, \dots, \quad (18)$$

where  $Y_k^* = Y^*(k)$  for  $k = 0, 1, 2, \dots$ , and  $\{Y^*(t), t \geq 0\}$  is a Brownian motion with mean

$$\theta_N = \Lambda - N\mu = -N\mu(1 - \rho) \quad (19)$$

and variance

$$\sigma_N^2 = \Lambda + \rho N\mu(1 - \mu) = \rho N\mu(2 - \mu). \quad (20)$$

Note that (19) and (20) are different from the mean  $-\mu\beta$  and variance  $\Lambda^* + \mu(1 - \mu)$  in Theorem 1, for two reasons: first, the process  $X_k^*$  and  $Y_k^*$  are diffusion *approximations* instead of the limiting processes stated in Theorem 1, which is why the term  $\rho$  appears in (19) and (20); second, the process  $X_k^*$  is to approximate the *centered* midnight count process (defined as  $\hat{X}_k = X_k - N$ ), not the diffusion-scaled version as in (2).

In the next three subsections, we first specify the basic adjoint relationship (BAR) for this discrete-time diffusion process  $X_k^*$ . Then, we discuss two ways to numerically calculate/approximate the stationary distribution of  $X_k^*$ : (i) a projection algorithm that numerically solves the BAR, and (ii) an approximate formula.

## 2.1 Basic adjoint relationship

The state space of  $\{X_k^*, k = 0, 1, \dots\}$  is  $\mathbb{R}$ . One can check that  $\{X_k^*, k = 0, 1, \dots\}$  is a discrete-time Markov process, since

$$X_{k+1}^* - X_k^* = Y_{k+1}^* - Y_k^* + \mu(X_k^*)^-, \text{ for } k = 0, 1, \dots,$$

and  $\{Y_{k+1}^* - Y_k^* : k = 0, 1, \dots\}$  is a sequence of iid normal r.v. with mean  $\theta_N$  and variance  $\sigma_N^2$ . The transition density of the Markov process is

$$p(x, y) = \mathbb{P}(X_{k+1}^* = y | X_k^* = x) = \begin{cases} \phi_{\theta_N, \sigma_N^2}(y - x), & \text{when } x \geq 0, \\ \phi_{\theta_N, \sigma_N^2}(y - (1 - \mu)x), & \text{when } x < 0, \end{cases} \quad (21)$$

where  $\phi_{\theta, \sigma^2}$  denotes the normal density function with mean  $\theta$  and variance  $\sigma^2$ . Let  $C_b(\mathbb{R})$  denote the set of bounded, continuous functions on  $\mathbb{R}$ . For each  $f \in C_b(\mathbb{R})$ , define

$$\mathbf{P}f(x) = \int_{\mathbb{R}} p(x, y) f(y) dy \quad \text{for each } x \in \mathbb{R}.$$

One can check that  $\mathbf{P}f \in C_b(\mathbb{R})$ . It follows that the stationary density  $\pi(x)$  satisfies

$$\int_{\mathbb{R}} \mathbf{P}f(x) \pi(x) dx = \int_{\mathbb{R}} f(x) \pi(x) dx, \quad \forall f \in C_b(\mathbb{R}), \quad (22)$$

or equivalently,

$$\int_{\mathbb{R}} \mathbf{L}f(x) \pi(x) dx = 0, \quad \forall f \in C_b(\mathbb{R}), \quad (23)$$

with  $\mathbf{L}f(x) = \mathbf{P}f(x) - f(x)$ . We call (23) the basic adjoint relationship (BAR) that governs the stationary density of the discrete-time Markov process  $\{X_k^*, k = 0, 1, \dots\}$ .

## 2.2 A projection algorithm

The BAR (23) is in the same format as (2.5) of [5]; the latter BAR is for the stationary density of a (continuous-time) diffusion process. As such the algorithm developed in [5] can be applied to compute the stationary density  $\pi^*$  of the discrete-time diffusion process  $\{X_k^*, k = 0, 1, \dots\}$ . We outline the algorithm here, commenting on the differences when appropriate.

### 2.2.1 Reference density and the space $L^2(\mathbb{R}, r)$

To compute the stationary density  $\pi^*$ , we first need a reference density  $r$  such that

$$\int_{\mathbb{R}} r(x) dx = 1.$$

We use the approximate formula  $\tilde{\pi}$  in Section 4.3.2 of the main paper (also see 34 below) as the reference density  $r$ .

Next, we define the ratio function as:

$$q(x) = \frac{\pi^*(x)}{r(x)} \quad \text{for } x \in \mathbb{R}. \quad (24)$$

With the given reference density  $r$ , if we can compute the ratio function  $q$ , then we can compute the stationary density via

$$\pi^*(x) = q(x)r(x) \text{ for } x \in \mathbb{R}.$$

To compute  $q$ , we plug (24) into (23) and get

$$\int_{\mathbb{R}} \mathbf{L}f(x)q(x)r(x)dx = 0, \quad \forall f \in C_b(\mathbb{R}). \quad (25)$$

Following the notation in [5], we use  $L^2(\mathbb{R}, r)$  to denote the space of all square-integrable functions on  $\mathbb{R}$  with respect to the measure that has density  $r$ . Namely,  $L^2(\mathbb{R}, r)$  is the set of measurable functions  $f$  on  $\mathbb{R}$  that satisfy

$$\int_{\mathbb{R}} f^2(x)r(x)dx < \infty.$$

We adopt the same inner product on  $L^2(\mathbb{R}, r)$  as in [5], that is,

$$\langle f, \hat{f} \rangle = \int_{\mathbb{R}} f(x)\hat{f}(x)r(x)dx, \quad \text{for } f, \hat{f} \in L^2(\mathbb{R}, r). \quad (26)$$

In (3.2) of [5], the authors made an important assumption on the reference density. Namely, they assumed that the reference density was chosen so that

$$q \in L^2(\mathbb{R}, r). \quad (27)$$

With our choice of the reference density  $r$ , we have been unable to verify that condition (27) is satisfied. We leave it as a conjecture that condition (27) is satisfied. The remainder of this section assumes that the conjecture is true.

### 2.2.2 Orthogonal projection

Note that the BAR (25) is equivalent to

$$\langle \mathbf{L}f, q \rangle = 0 \quad \text{for each } f \in C_b(\mathbb{R}).$$

Thus,  $q$  satisfying the BAR is equivalent to  $q$  being *orthogonal* to  $\mathbf{L}f$  for each  $f \in C_b(\mathbb{R})$ . We define a space  $H$  as

$$H = \text{the closure of } \{\mathbf{L}f : f \in C_b(\mathbb{R})\},$$

which is a subspace of  $L^2(\mathbb{R}, r)$ . Therefore,  $q$  satisfying the BAR is equivalent to  $q$  being orthogonal to space  $H$ . Therefore, our task is to find a function  $q$  that is orthogonal to space  $H$ . To do so, we consider a constant function  $e$  with  $e(x) = 1$  for each  $x \in \mathbb{R}$ . Since

$$\langle e, q \rangle = \int_{\mathbb{R}} e(x)q(x)r(x)dx = \int_{\mathbb{R}} \pi^*(x)dx = 1, \quad (28)$$

one can check that  $e \notin H$  because otherwise  $\langle e, q \rangle = 0$ , contradicting (28). We use  $\bar{e}$  to denote the projection of  $e$  onto  $H$ . Then,  $e - \bar{e} \neq 0$  and it must be orthogonal to  $H$ . Once we have  $\bar{e}$ , we obtain the ratio function  $q$  by

$$q = \frac{e - \bar{e}}{\|e - \bar{e}\|^2},$$

where  $\|\cdot\|$  is the induced norm from the inner product (26) with  $\|f\|^2 = \langle f, f \rangle$  for  $f \in L^2(\mathbb{R}, r)$ .

### 2.2.3 Finite-dimensional approximation

The projection of  $e$  onto  $H$  can be expressed as

$$\bar{e} = \operatorname{argmin}_{h \in H} \|e - h\|. \quad (29)$$

The space  $H$  is linear and infinitely dimensional. To compute the projection numerically, we use a finite-dimensional subspace  $H_k$  to approximate  $H$  and find the projection  $\bar{e}_k$  of  $e$  on  $H_k$ , namely,

$$\bar{e}_k = \operatorname{argmin}_{h \in H_k} \|e - h\|. \quad (30)$$

Let  $C_k$  be a finite-dimensional, linear subspace of  $C_b(\mathbb{R})$ . Then  $H_k = \{\mathbf{L}f : f \in C_k\}$  is a finite-dimensional subspace of  $H$ . Assume that  $\{f_i : i = 1, 2, \dots, m\} \subset C_k$  is a basis of  $C_k$ . Then, since the projection  $\bar{e}_k \in H_k$ , it can be represented as a linear combination of  $\{\mathbf{L}f_i : i = 1, 2, \dots, m\}$ . That is,

$$\bar{e}_k = \sum_{i=1}^m \alpha_i \mathbf{L}f_i \quad (31)$$

where  $\alpha_i \in \mathbb{R}$  for  $i = 1, 2, \dots, m$ .

To compute the vector of coefficients  $\alpha = (\alpha_1, \dots, \alpha_m)'$ , we use the fact that  $\langle e - \bar{e}_k, \mathbf{L}f_i \rangle = 0$  for  $i = 1, 2, \dots, m$ . Consequently, we obtain a system of linear equations

$$A\alpha = \beta, \quad (32)$$

where  $A_{ij} = \langle \mathbf{L}f_i, \mathbf{L}f_j \rangle$  and  $\beta_i = \langle e, \mathbf{L}f_i \rangle$  for  $i, j = 1, \dots, m$ . The matrix is symmetric, semi-positive definite, but can be singular. Although the solution to the system of linear equations may not be unique, projection  $\bar{e}_k$  is unique. When  $A$  is singular or nearly singular, one can solve (32) by direct methods such as the QR decomposition and the Cholesky decomposition or by iterative methods such as LSQR [10]. The Cholesky decomposition exploits the symmetric and semi-positive definite properties of  $A$  even when  $A$  is singular [1, 8], whereas QR decomposition does not. Unlike many other iterative methods, LSQR can handle matrix  $A$  when it is singular. LSQR does not exploit semi-positive definiteness.

Once we get the vector of coefficients  $\alpha = (\alpha_1, \dots, \alpha_m)'$  by solving the system of linear equations (32), we can compute  $\bar{e}_k$  as in (31). Eventually, we can approximately compute the stationary density  $\pi^*$  as

$$\pi^*(x) \approx r(x) \frac{1 - \bar{e}_k(x)}{\|e - \bar{e}_k\|^2} \quad \forall x \in \mathbb{R}. \quad (33)$$

### 2.2.4 FEM implementation

In our implementation, we use the finite element method (FEM) to construct the approximate space  $C_k$ , following Section 3.3 of [5]. The numerical results in this paper for approximating the stationary density  $\pi$  with the projection algorithm all follow this FEM implementation. In Proposition 3 of Dai and He [5], they proved the convergence of using (33) to approximate  $\pi$  as  $H_k \uparrow H$ . Their proof applies to our setting when (27) is satisfied.

## 2.3 Approximate formula for the stationary density

In Section 4.3 of the main paper, the following formula  $\tilde{\pi}$  is proposed as a proxy for the stationary density  $\pi^*$  of the diffusion process  $X^*$ :

$$\tilde{\pi}(x) = \begin{cases} \alpha_1 \exp(2\theta_N x / \sigma_N^2), & x \geq 0; \\ \alpha_2 \exp(-(2\mu - \mu^2)(x - \theta_N / \mu)^2 / 2\sigma_N^2), & x < 0; \end{cases} \quad (34)$$



where  $\alpha_1$  and  $\alpha_2$  are normalizing constants that make  $\tilde{\pi}(x)$  continuous at zero and  $\int_{\mathbb{R}} \tilde{\pi}(x)dx = 1$ .

As mentioned in the main paper, the rationale of this approximate formula is based on the analogy between  $\{X_k^* : k = 0, 1, 2, \dots\}$  and  $\{\tilde{X}(t), t \geq 0\}$ , where

$$\tilde{X}(t) = \tilde{Y}(t) + \mu \int_0^t (\tilde{X}(s))^- ds, \quad t \geq 0, \quad (35)$$

and  $\{\tilde{Y}(t), t \geq 0\}$  is a Brownian motion. To get the stationary density of  $\tilde{X}$ , Browne and Whitt [3] have suggested that since (i)  $\tilde{X}$  is a Ornstein-Uhlenbeck (OU) process on  $(-\infty, 0]$  and the stationary density of an OU process has a Gaussian form and (ii)  $\tilde{X}$  is a reflected Brownian motion (RBM) on  $[0, \infty)$  and the stationary density of a RBM has an exponential form, then the stationary density of  $\tilde{X}$  can be obtained by piecing together the Gaussian and exponential densities.

We use the same piecing technique in our setting. Specifically,  $\{X_k^* : k = 0, 1, 2, \dots\}$  behaves as a discrete version of the OU process on  $(-\infty, 0]$  and as a reflected random walk on  $[0, \infty)$ . We show in Proposition 1 below that the stationary density of the discrete-time OU (DOU) process also has a Gaussian form. For the reflected random walk, existing research shows that it has an exponential tail [9, 11, 2]. Therefore, we piece together a Gaussian density and an exponential density and propose using (34) to approximate  $\pi^*$ . In the next two subsections, we first prove that the stationary density of the discrete-time OU process has a Gaussian form, then we show the details of deriving formula (34).

### 2.3.1 The stationary distribution of a discrete OU process

Similar to the continuous-time version of the Ornstein-Uhlenbeck process, we define its discrete-time version  $\{X_k^{\text{DOU}}, k = 0, 1, \dots\}$  as:

$$X_k^{\text{DOU}} = Y_k^{\text{DOU}} - \mu \sum_{i=0}^{k-1} X_i^{\text{DOU}}, \quad k = 0, 1, \dots \quad (36)$$

where  $\{Y_k^{\text{DOU}} := \sum_{i=0}^{k-1} \xi_i, k = 0, 1, \dots\}$  is a Gaussian random walk, i.e.,  $\{\xi_i\}$  is a sequence of iid random variables following a normal distribution with mean  $\theta$  and variance  $\sigma^2$ .

The following proposition says the stationary density for a discrete OU process has the Gaussian form, which is consistent with that in a continuous-time OU process.

**Proposition 1** *Given  $0 < \mu < 1$ , for a discrete-time Ornstein-Uhlenbeck (DOU) process  $\{X_k^{\text{DOU}}, k = 0, 1, \dots\}$  satisfying*

$$X_k^{\text{DOU}} = Y_k^{\text{DOU}} - \mu \sum_{i=0}^{k-1} X_i^{\text{DOU}}, \quad k = 0, 1, \dots \quad (37)$$

*where  $\{Y_k^{\text{DOU}}\}$  is a Gaussian random walk with drift  $\theta$  and variance  $\sigma^2$ , the stationary density of the DOU process,  $\pi$ , is a normal density with mean  $\theta/\mu$  and variance  $\frac{\sigma^2}{2\mu - \mu^2}$ .*

*Proof for Proposition 1.* Note that the DOU process  $\{X_k^{\text{DOU}}, k = 0, 1, \dots\}$  satisfying (37) is a Markov process since

$$X_{k+1}^{\text{DOU}} - X_k^{\text{DOU}} = (Y_{k+1}^{\text{DOU}} - Y_k^{\text{DOU}}) - \mu X_k^{\text{DOU}}.$$

The transition probability from state  $y$  to state  $x$  is

$$\mathbb{P}(X_{k+1}^{\text{DOU}} = x | X_k^{\text{DOU}} = y) = \phi_{\theta, \sigma^2}(x - (1 - \mu)y),$$

where  $\phi_{\theta, \sigma^2}(s)$  denotes the probability density function associated with a normal random variable with mean  $\theta$  and variance  $\sigma^2$ .

To prove this proposition, we just need to show that

$$\pi(x) = \int_{-\infty}^{\infty} \mathbb{P}(x|y)\pi(y)dy \quad (38)$$

for any given  $x$ , where

$$\pi(x) = \frac{\sqrt{(2\mu - \mu^2)}}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(2\mu - \mu^2)(x - \theta/\mu)^2}{2\sigma^2}\right).$$

We have

$$\begin{aligned} \mathbb{P}(x|y)\pi(y) &= \frac{\sqrt{2\mu - \mu^2}}{\sqrt{2\pi}\sigma} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(2\mu - \mu^2)(y - \theta/\mu)^2}{2\sigma^2}\right) \exp\left(-\frac{(x - (1 - \mu)y - \theta)^2}{2\sigma^2}\right) \\ &= \frac{\sqrt{(2\mu - \mu^2)}}{\sqrt{2\pi}\sigma} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{y^2 - 2[(1 - \mu)x + \theta]y}{2\sigma^2}\right) \exp\left(-\frac{(2\mu - \mu^2)\theta^2/\mu^2 + (x - \theta)^2}{2\sigma^2}\right) \\ &= \frac{\sqrt{(2\mu - \mu^2)}}{\sqrt{2\pi}\sigma} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{[y - ((1 - \mu)x + \theta)]^2}{2\sigma^2}\right) \\ &\quad \cdot \exp\left(-\frac{(2\mu - \mu^2)\theta^2/\mu^2 + (x - \theta)^2 - [(1 - \mu)x + \theta]^2}{2\sigma^2}\right). \end{aligned}$$

Among which,

$$\begin{aligned} V(x) &= \exp\left(-\frac{(2\mu - \mu^2)\theta^2/\mu^2 + (x - \theta)^2 - [(1 - \mu)x + \theta]^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{(2\mu - \mu^2)\theta^2/\mu^2 + (2\mu - \mu^2)x^2 - 2(2 - \mu)\theta x}{2\sigma^2}\right) \\ &= \exp\left(-\frac{(2\mu - \mu^2)(x - \theta/\mu)^2}{2\sigma^2}\right). \end{aligned}$$

Then, we have

$$\begin{aligned} \int_{-\infty}^{\infty} \mathbb{P}(x|y)\pi(y)dy &= \frac{\sqrt{(2\mu - \mu^2)}}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(2\mu - \mu^2)(x - \theta/\mu)^2}{2\sigma^2}\right) \\ &\quad \cdot \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{[y - ((1 - \mu)x + \theta)]^2}{2\sigma^2}\right) dy \\ &= \frac{\sqrt{(2\mu - \mu^2)}}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(2\mu - \mu^2)(x - \theta/\mu)^2}{2\sigma^2}\right), \end{aligned}$$

which takes the exact form as the normal density with mean  $\theta/\mu$  and variance  $\sigma^2/(2\mu - \mu^2)$  and thus, equals to  $\pi(x)$ . This completes our proof for  $\pi$  being the stationary density.

### 2.3.2 Derivation of the approximate formula

Based on Proposition 1, we conjecture that the stationary distribution of  $X^*$  can be approximated by the following form:

$$\tilde{\pi}(x) = \begin{cases} \pi_1(x) = \alpha_1 \exp(-\gamma x), & x \geq 0; \\ \pi_2(x) = \alpha_2 \exp(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}), & x < 0. \end{cases} \quad (39)$$

For the ease of exposition, we use  $\theta$  and  $\sigma^2$  instead of  $\theta_N$  and  $\sigma_N^2$  to denote the mean and variance of the discrete-time diffusion process  $X^*$ . Moreover, in (39),  $\alpha_1$  and  $\alpha_2$  are two normalizing constants,  $\gamma$  is the unknown parameter for the exponential density part, and we define

$$\beta = -\mu/\theta.$$

The stationary density should satisfy

$$\tilde{\pi}(x) = \int_{-\infty}^{\infty} p(y, x) \tilde{\pi}(y) dy,$$

one special form of the BAR (22), or equivalently,

$$\tilde{\pi}(x) = \int_0^{\infty} p(y, x) \pi_1(y) dy + \int_{-\infty}^0 p(y, x) \pi_2(y) dy, \quad (40)$$

where  $p(y, x)$  is the transitional density of  $X^*$  (from state  $y$  to state  $x$ ) defined in (21).

We rewrite Equation(40) as follows. First, for  $y \geq 0$ , we have

$$\begin{aligned} p(y, x) \pi_1(y) &= \frac{\alpha_1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(x - y + \mu\beta)^2}{2\sigma^2}\right) \cdot \exp(-\gamma y) \\ &= \frac{\alpha_1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{y^2 - 2(x + \mu\beta - \sigma^2\gamma)y + (x + \mu\beta)^2}{2\sigma^2}\right) \\ &= \frac{\alpha_1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{[y - (x + \mu\beta - \sigma^2\gamma)]^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{(x + \mu\beta)^2 - (x + \mu\beta - \sigma^2\gamma)^2}{2\sigma^2}\right) \\ &= \frac{\alpha_1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{[y - (x + \mu\beta - \sigma^2\gamma)]^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{\sigma^2\gamma[2x + (2\mu\beta - \sigma^2\gamma)]}{2\sigma^2}\right) \\ &= \frac{\alpha_1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) \cdot \exp\left(-\frac{[y - (x + \mu\beta - \sigma^2\gamma)]^2}{2\sigma^2}\right) \cdot \exp(-\gamma x). \end{aligned}$$

Therefore,

$$\begin{aligned} \int_0^{\infty} p(y, x) \pi_1(y) dy &= \alpha_1 \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) \exp(-\gamma x) \cdot \int_0^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{[y - (x + \mu\beta - \sigma^2\gamma)]^2}{2\sigma^2}\right) dy \\ &= \alpha_1 \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) \exp(-\gamma x) \cdot [1 - \Phi_{-\mu\beta, \sigma^2}(-x - (2\mu\beta - \sigma^2\gamma))]. \end{aligned}$$

Second, for  $y < 0$ , we have

$$\begin{aligned} p(y, x) \pi_2(y) &= \frac{\alpha_2}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(x - (1 - \mu)y + \mu\beta)^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{(2\mu - \mu^2)(y + \beta)^2}{2\sigma^2}\right) \\ &= \frac{\alpha_2}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(y - ((1 - \mu)x - \mu\beta))^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right). \end{aligned}$$

Therefore,

$$\begin{aligned} \int_{-\infty}^0 p(y, x) \pi_2(y) dy &= \int_{-\infty}^0 \frac{\alpha_2}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{(y - ((1 - \mu)x - \mu\beta))^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right) dy \\ &= \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right) \cdot \int_{-\infty}^0 \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - (1 - \mu)x + \mu\beta)^2}{2\sigma^2}\right) dy \\ &= \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right) \cdot \Phi_{-\mu\beta, \sigma^2}(-x + \mu x). \end{aligned}$$

If (40) holds, when  $x \geq 0$ , we should have

$$\begin{aligned}\alpha_1 \exp(-\gamma x) &= \alpha_1 \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) \exp(-\gamma x) \cdot [1 - \Phi_{-\mu\beta, \sigma^2}(-x - (2\mu\beta - \sigma^2\gamma))] \\ &\quad + \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right) \cdot \Phi_{-\mu\beta, \sigma^2}(-x + \mu x),\end{aligned}$$

which is equivalent to

$$\begin{aligned}&\alpha_1 \exp(-\gamma x) \cdot \left[1 - \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) (1 - \Phi_{-\mu\beta, \sigma^2}(-x - (2\mu\beta - \sigma^2\gamma)))\right] \\ &= \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right) \cdot \Phi_{-\mu\beta, \sigma^2}(-x + \mu x).\end{aligned}\tag{41}$$

Similarly, if (40) holds, when  $x < 0$ , we should have

$$\begin{aligned}\alpha_2 \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right) &= \alpha_1 \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) \exp(-\gamma x) \cdot [1 - \Phi_{-\mu\beta, \sigma^2}(-x - (2\mu\beta - \sigma^2\gamma))] \\ &\quad + \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right) \cdot \Phi_{-\mu\beta, \sigma^2}(-x + \mu x),\end{aligned}$$

which is equivalent to

$$\begin{aligned}&\alpha_1 \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) \exp(-\gamma x) \cdot [1 - \Phi_{-\mu\beta, \sigma^2}(-x - (2\mu\beta - \sigma^2\gamma))] \\ &= \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)(x + \beta)^2}{2\sigma^2}\right) \cdot (1 - \Phi_{-\mu\beta, \sigma^2}(-x + \mu x)).\end{aligned}\tag{42}$$

When  $x = 0$ , Equations (41) and (42) become

$$\alpha_1 \cdot \left[1 - \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) (1 - \Phi_{-\mu\beta, \sigma^2}(-(2\mu\beta - \sigma^2\gamma)))\right] = \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)\beta^2}{2\sigma^2}\right) \cdot \Phi_{-\mu\beta, \sigma^2}(0),\tag{43}$$

and

$$\alpha_1 \exp\left(-\frac{\gamma(2\mu\beta - \sigma^2\gamma)}{2}\right) \cdot [1 - \Phi_{-\mu\beta, \sigma^2}(-(2\mu\beta - \sigma^2\gamma))] = \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)\beta^2}{2\sigma^2}\right) \cdot (1 - \Phi_{-\mu\beta, \sigma^2}(0)),\tag{44}$$

respectively.

Recall that  $\pi(x)$  is continuous at  $x = 0$ . Thus, the two normalizing constants satisfy:

$$\pi_1(0) = \alpha_1 = \alpha_2 \exp\left(-\frac{(2\mu - \mu^2)\beta^2}{2\sigma^2}\right) = \pi_2(0).\tag{45}$$

Comparing (45) with (43) and (44), we find that when

$$\sigma^2\gamma = 2\mu\beta = -2\theta,$$

or equivalently,

$$\gamma = -\frac{2\theta}{\sigma^2},\tag{46}$$

both (43) and (44) can be satisfied. Therefore, we choose  $\gamma$  in (46), which eventually gives us (34).

Unfortunately, using this  $\gamma$ , we are unable to show (41) and (42) hold for a general  $x$ .

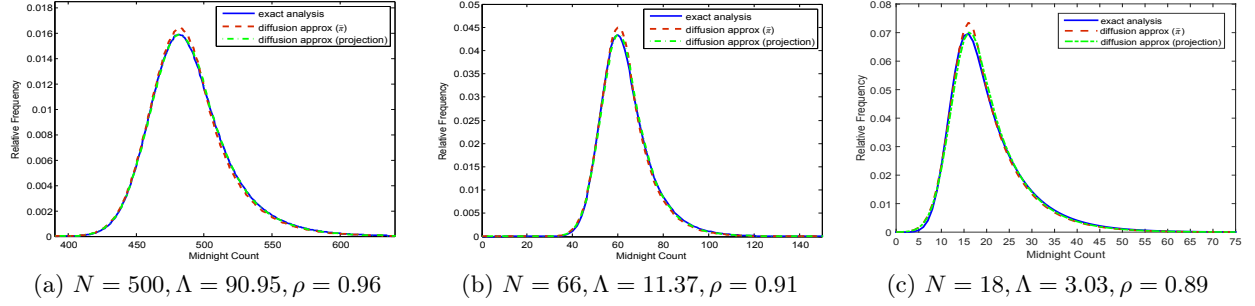


Figure 1: Stationary distribution of the midnight customer count from exact Markov chain analysis and diffusion approximations. Here, the mean LOS is 5.3 days, and we do not specify the discharge distribution because it does not affect the midnight customer count distribution.

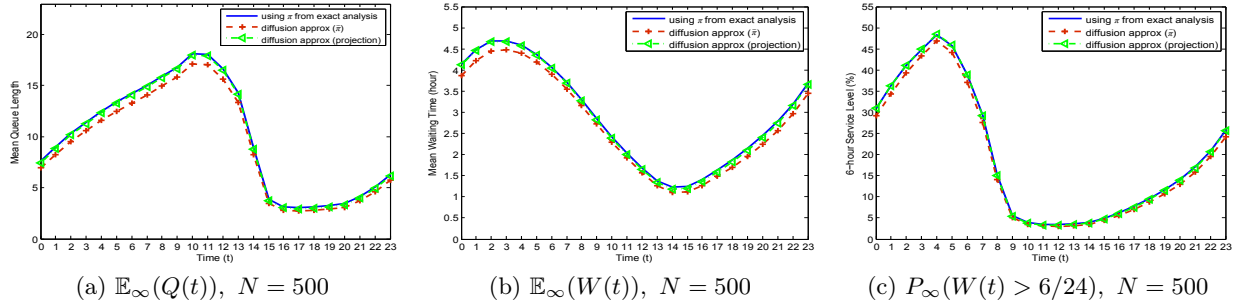


Figure 2: Time-dependent performance curves from exact analysis and diffusion approximations. Here,  $\Lambda = 90.95$  for  $N = 500$ . We fix the mean LOS as 5.3 days and use the baseline discharge distribution. The three performance curves in each subfigure are from normal approximations using (i)  $\pi$  solved from exact Markov chain analysis, (ii)  $\tilde{\pi}$  in (34), and (iii)  $\pi^*$  solved from the projection algorithm, respectively.

### 3 Numerical results on diffusion approximations

#### 3.1 Approximation for the midnight count distribution

Figure 1 compares the stationary distributions of the midnight customer count solved (i) from the exact Markov chain analysis, (ii) from using the approximate formula  $\tilde{\pi}$  in (34), and (iii) from using the projection algorithm specified in Section 2.2. The parameter settings for these numerical experiments are the same as those in Section 5 of the main paper. We test a large system ( $N = 500$ ) and two small systems ( $N = 66$  and  $18$ ), with the utilization  $\rho$  being 96%, 91% and 89%, respectively.

#### 3.2 Time-dependent performance

Figures 2, 3 and 4 show the time-dependent performance for systems with  $N = 500$ ,  $66$ , and  $18$ , respectively. The three curves in each subfigure are obtained from normal approximations using (i)  $\pi$  solved from exact Markov chain analysis, (ii)  $\tilde{\pi}$  in (34), and (iii)  $\pi^*$  solved from the projection algorithm specified in Section 2.2.

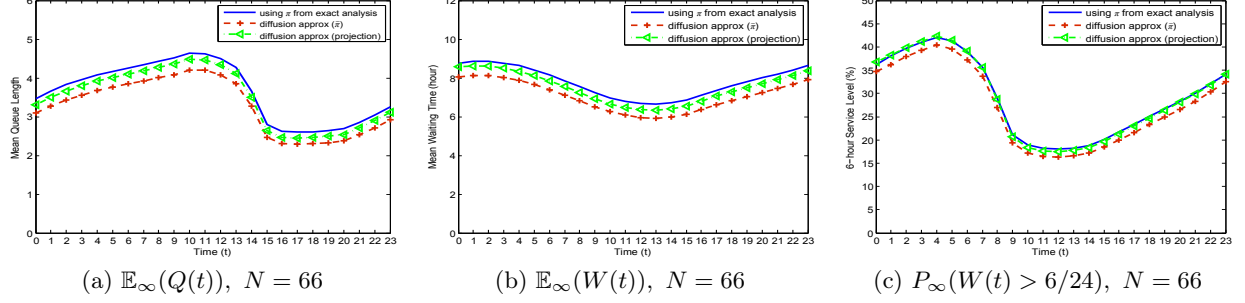


Figure 3: Time-dependent performance curves from exact analysis and diffusion approximations. Here,  $\Lambda = 11.37$  for  $N = 66$ . We fix the mean LOS as 5.3 days and use the baseline discharge distribution. The three performance curves in each subfigure are from normal approximations using (i)  $\pi$  solved from exact Markov chain analysis, (ii)  $\tilde{\pi}$  in (34), and (iii)  $\pi^*$  solved from the projection algorithm, respectively.

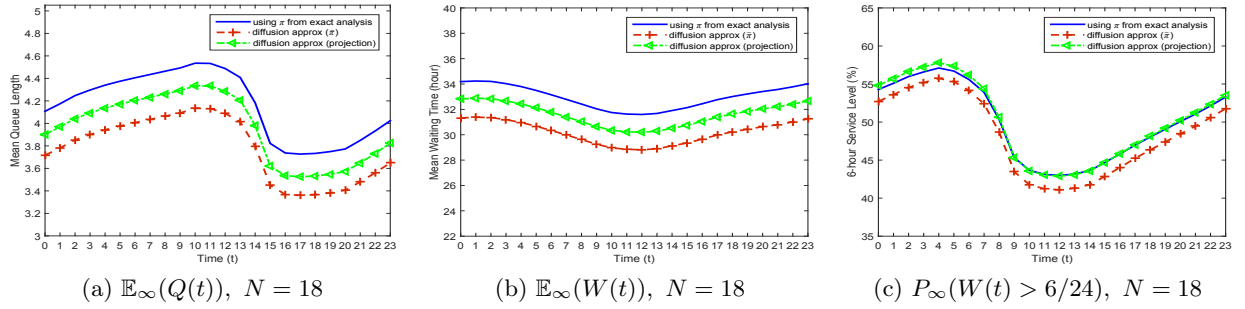


Figure 4: Time-dependent performance curves from exact analysis and diffusion approximations. Here,  $\Lambda = 3.03$  for  $N = 18$ . We fix the mean LOS as 5.3 days and use the baseline discharge distribution. The three performance curves in each subfigure are from normal approximations using (i)  $\pi$  solved from exact Markov chain analysis, (ii)  $\tilde{\pi}$  in (34), and (iii)  $\pi^*$  solved from the projection algorithm, respectively.

## References

- [1] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammerling, A. McKenney *et al.*, *LAPACK Users' guide*. SIAM, 1999, vol. 9.
- [2] J. Blanchet and P. Glynn., "Complete corrected diffusion for the maximum of the random walk." *Annals of Applied Probability*, vol. 16, pp. 951–953, 2006.
- [3] S. Browne and W. Whitt, "Piecewise-linear diffusion processes," in *Advances in Queueing*, J. Dshalalow, Ed. Boca Raton, FL: CRC Press, 1995, pp. 463–480.
- [4] J. G. Dai, S. He, and T. Tezcan, "Many-server diffusion limits for  $G/Ph/n + GI$  queues," *Annals of Applied Probability*, vol. 20, no. 5, pp. 1854–1890, 2010.
- [5] J. G. Dai and S. He, "Many-server queues with customer abandonment: Numerical analysis of their diffusion model," *Stochastic Systems*, vol. 3, no. 1, pp. 96–146, 2013.
- [6] J. G. Dai and P. Shi, "A two-time-scale approach to time-varying queues for hospital inpatient flow management," 2014, working paper. [Online]. Available: [http://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2489533](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2489533)
- [7] N. Gans, G. Koole, and A. Mandelbaum, "Telephone call centers: Tutorial, review, and research prospects," *Manufacturing & Service Operations Management*, vol. 5, no. 2, pp. 79–141, 2003.
- [8] S. Hammarling, N. J. Higham, C. Lucas, M. Eprint, S. Hammarling, N. J. Higham, and C. Lucas, "LAPACK-style codes for pivoted Cholesky," 2007.
- [9] J. F. C. Kingman, "Ergodic properties of continuous-time markov processes and their discrete skeletons," *Proceedings of the London Mathematical Society*, vol. s3-13, no. 1, pp. 593–604, 1963.
- [10] C. C. Paige and M. A. Saunders, "Lsqqr: An algorithm for sparse linear equations and sparse least squares," *ACM Transactions on Mathematical Software (TOMS)*, vol. 8, no. 1, pp. 43–71, 1982.
- [11] D. Siegmund, "Corrected diffusion approximations in certain random walk problems," *Advances in Applied Probability*, vol. 11, no. 4, pp. 701–719, 1979.